

Deep Learning for Bias Detection: From Inception to Deployment

Md Abul Bashar¹[0000–0003–1004–4085], Richi Nayak¹, Anjor Kothare², Vishal Sharma², and Kesavan Kandadai²

¹ Queensland University of Technology, Brisbane, Australia
{m1.bashar, r.nayak}@qut.edu.au
² ishield.ai

anjor.kothare@ishield.ai, vishal.sharma@tothenew.com,
kesavan@pallavatech.com, kesavan@ishield.ai

Abstract. To create a more inclusive workplace, enterprises are actively investing in identifying and eliminating unconscious bias (e.g., gender, race, age, disability, elitism and religion) across their various functions. We propose a deep learning model with a transfer learning based language model to learn from manually tagged documents for automatically identifying bias in enterprise content. We first pretrain a deep learning-based language-model using Wikipedia, then fine tune the model with a large unlabelled data set related with various types of enterprise content. Finally, a linear layer followed by softmax layer is added at the end of the language model and the model is trained on a labelled bias dataset consisting of enterprise content. The trained model is thoroughly evaluated on independent datasets to ensure a general application. We present the proposed method and its deployment detail in a real-world application.

Keywords: Deep Learning · Bias Detection · Transfer Learning · Text Data · Small Dataset.

1 Introduction

A rigorous study by McKinsey [19] found that more diverse and inclusive organisations outperform those that are not. The concept of unconscious bias has become increasingly pervasive, with many organisations training their employees on the concept of Diversity and Inclusion (D&I). Unfortunately, though an increased awareness of unconscious bias can have benefits, it is not a systemic and consistent solution [25]. Enterprises are looking for technologies that can create consistent and scalable practices to identify or mitigate biases across organisations, often in real-time.

While there have been a few analytics products to measure employee demographics, pay parity and customised training [40], there is lack of solutions targeted toward tackling unconscious bias in content that is created within an enterprise. For every enterprise, almost always the first touch point for their users, vendors, employees or business partners is the content published across their digital channels. This content may be job descriptions, website content,

marketing messages, blogs, reports and social media content. Internally within the enterprise, content is also created in the form of product documentation, user guides, customer care, messaging and collaboration platforms. Unconscious bias in content manifests itself in various forms in different types of content. For example: (a) in job descriptions this could be through use of pronouns or masculine coded words, (b) in enterprise messaging platforms this could be through toxic messages and bullying, (c) in product descriptions this could be through stereotyping a type of a user, (d) in e-commerce product details this could be through references to body shapes, or (e) in a stock market report, this could be through inherent assumption of the gender of a potential investor.

In this paper, we aim to propose a Deep Learning (DL) model to detect unconscious bias in content and how it can be applied to integrate and work with enterprise applications. Lexical based systems and traditional machine learning systems do not solve the problem due to complexity and intricate relationships inherent in the textual narratives. With advancements in DL in Natural Language Processing (NLP), it is feasible to build a DL based system. However, such a system requires a large set of labelled data for building the model. With the manual efforts involved in labelling the data, it is difficult to create a large dataset. There do not exist similar data that can be used in labelling. In this paper we propose a novel method based on transfer learning to deal with the small set of labelled data and detect unconscious bias in content.

This research makes the following novel contributions. (1) It proposes a comprehensive bias detection model that can detect four types of bias (*Race, Gender, Age, Not Appropriate*) in text data. (2) It uses progressive transfer learning that allows to train a smaller model (i.e. less number of hyper parameters) with a small labelled dataset for better accuracy. (3). It presents the detail of deploying the model in a real-world application³.

2 Literature review

When the unconscious biases are not tackled, they cause serious harm to the business including financially, socially and culturally [1,47,6,19,29,3]. However, very few to no technology solutions exist to tackle unconscious bias in the content that is created and published across the enterprise for both their internal and external audiences [21].

Deep learning models have become quite successful in natural language processing, e.g., content generation, language translation, Question and Answering systems, text classification [8,7,10,9,2] and clustering. In spite of this success of DL in NLP, there has been very limited works on building a generic bias detection method. The work in current literature can be grouped into two categories: (1) text representation learning; and (2) reducing subjective bias in text.

Researchers⁴ [34,12] proposed to debias semantic representation of words by removing bias component from word embedding. The goal is to make semantic representations fair across attributes like gender and age. Autoencoder was used to generate a balanced gender-oriented word distribution to remove gender bias

³ <https://ishield.ai/>

⁴ <https://gender-decoder.katmatfield.com/>

from word embeddings[31]. Counterfactual data has been augmented to alter the training distribution to balance gender-based statistics [35,48,17]. Research in [45] used adversarial training to squeeze directions of bias in the hidden states of image representation.

Research in [42] proposed a model for automatically suggesting edits for subjective-bias words following the Wikipedia’s neutral point of view (NPOV) policy for defining subjective bias. However, the model is limited to single-word edits, i.e. it can handle simplest instances of bias only. Research in [28,43] used lexicon of bias words for detecting language bias in sentences of Wikipedia articles. A fine-tuned BERT model [16] is used for detecting gender bias only [18].

An issue faced by bias detection methods in multi-class setting is failing to select minority class examples in imbalanced data distribution [5]. Guided learning, based on crowd-sourcing to find or generate class-specific training instances, can help to get more balanced class frequencies [41,4]. However, guided learning is resource consuming and may not present the true distribution generating training examples. Sometimes, heuristic labelling methods such as distant supervision [15] or data programming [20] are used for datasets with the imbalanced class distribution. However, these methods are only applicable when a good knowledge base or a pretrained predictor is available [14].

In this research, we propose to use progressive transfer learning to address the class imbalance problem. To our best of knowledge, there exist no model that comprehensively detect common biases (e.g. race, gender, age) in text.

3 Language based Deep Learning model for Bias Detection

Bias detection in the text data is a complex problem because usually bias is represented by linguistic cues that are subtle and can be determined only through its context in the text. Let X be a text dataset that contains n features and K classes. Let $\mathbf{x} = \langle x_1, \dots, x_n \rangle$ be a vector representing an instance in X . Let C_k be a set of K classes. The bias detection is a classification task that assigns an instance to a bias class (or category) C_k based on the feature vector \mathbf{x} ; i.e. $f \in \mathcal{F} : X \rightarrow C_k$, where $f(\mathbf{x}) = \max_{C_k} p(C_k|\mathbf{x})$. This ascertains that we need to know $p(C_k|\mathbf{x})$ for a bias detection task. The joint probability $p(\mathbf{x}, C_k)$ of \mathbf{x} and C_k can be written as

$$p(\mathbf{x}, C_k) = p(C_k|\mathbf{x})p(\mathbf{x}) \tag{1}$$

where $p(\mathbf{x})$ is the prior probability distribution. The prior probability $p(\mathbf{x})$ can be seen as a regulariser for $p(C_k|\mathbf{x})$ that can regularise modelling of the associated uncertainties of $p(\mathbf{x}, C_k)$ [10]. As $p(\mathbf{x})$ does not depend on C_k , this means that $p(\mathbf{x})$ can be learned independent of the class level C_k . That is, $p(\mathbf{x})$ can be learned from unlabelled data.

Prior research [10,9] showed that the estimation of $p(\mathbf{x}, C_k)$ can be improved when $p(\mathbf{x})$ is learned from a sequence of unlabelled datasets, especially when the labelled dataset is small. In this study, we propose to implement a bias classification model utilising this technique of improving the prediction accuracy through unlabelled data.

A discriminative model such as LSTM learns to classify an instance \mathbf{x} into class C_k by learning the conditional probability distribution as $p(C_k|\mathbf{x}, \theta) \approx p(C_k|\mathbf{x})$ where θ is the list of model parameters. However, accurately approximating $p(C_k|\mathbf{x})$ requires a large number of labelled instances. If only a small set of labelled data is available, the learned $p(C_k|\mathbf{x}, \theta)$ might not be a good approximation of population distribution because θ may over-fit the small training set. Alternately, $p(\mathbf{x})$ can be learned from one or more large unlabelled datasets and conditioned on \mathbf{x} to learn $p(C_k|\mathbf{x}, \theta)$, leading to $p(C_k|\mathbf{x}, \theta)p(\mathbf{x}) \approx p(\mathbf{x}, C_k)$. The term $p(C_k|\mathbf{x}, \theta)p(\mathbf{x})$ can be seen as equivalent to combining the regularisation into the discriminative model. This regularised model would act similar to a generative model.

Unlike common transfer learning where $p(\mathbf{x})$ is learned once from a large unlabelled dataset, we propose to progressively learn $p(\mathbf{x})$ from a sequence of unlabelled datasets. This allows us to use a relatively small set of parameters in our model. Then we use $p(\mathbf{x})$ to learn $p(C_k|\mathbf{x}, \theta)p(\mathbf{x}) \approx p(\mathbf{x}, C_k)$ with a small training dataset. Next, we present the estimation of $p(\mathbf{x})$ as a neural network language model (NNLM) using unsupervised learning.

3.1 Neural Network Language Model

Probability $p(\mathbf{x})$ can be estimated using the assumption of Language model where features are considered conditionally dependent [9,10]. This is to support natural language processing where in a sentence, the sequencing of words depends on each other. Based on this, the joint probability $p(\mathbf{x}, C_k)$ in Equation 1 can be rewritten as follows, using the chain rule:

$$\begin{aligned} p(\mathbf{x}, C_k) &= p(C_k|\mathbf{x})p(\mathbf{x}) \\ &= p(C_k|\mathbf{x})p(x_1, \dots, x_n) \\ &= p(C_k|\mathbf{x}) \prod_{i=1}^n p(x_i|x_1 \dots x_{i-1}) \end{aligned} \tag{2}$$

The part $\prod_{i=1}^n p(x_i|x_1 \dots x_{i-1})$ in Equation 2 can be considered as a language model because $p(x_i|x_1 \dots x_{i-1})$ seeks to predict the probability of observing the i th feature x_i , given the previous $(i - 1)$ features $(x_1 \dots x_{i-1})$. A Recurrent Neural Network (RNN) or its variants such as Long Short-Term Memory (LSTM) can be used to model $\prod_{i=1}^n p(x_i|x_1 \dots x_{i-1})$ as they work in a similar way to capture the order of features and their non-linear and hierarchical interactions [39,30,10]. Similar to traditional language models, a RNN/LSTM based NNLM can approximate joint probabilities over the feature sequences as follows, where ω is the list of model parameters.

$$\begin{aligned} p(\mathbf{x}) &= p(x_1, \dots, x_n) \\ &\approx p(x_1, \dots, x_n, \omega) \\ &\approx \prod_{i=1}^n p(x_i|x_1 \dots x_{i-1}, \omega) \end{aligned} \tag{3}$$

Given a sequence of features, a RNN recurrently processes each feature and uses multiple hidden layers to capture the order of features and their non-linear and hierarchical interactions. The hidden state is used to derive a vector of probabilities representing the network’s guess of the subsequent feature in the sequence [10,9]. The network aims to minimize the loss calculated based on the vector of probabilities and the actual next feature. In simple words, the context of all previous features in the sequence is encoded within the parameters ω of the network and the probability of getting the next word is distributed over the vocabulary using a Softmax function [30].

Estimating $p(\mathbf{x})$ using a LSTM-based LM (i.e. NNLM) model on a huge dataset that covers multitude of domains is useful for transfer learning, but it can be very expensive in terms of required computation and memory [13]. Additionally, it can learn irrelevant and misleading relationships in data due to interactions between different domains in a single corpus [10]. Therefore, as suggested in [9], we use an alternative approach based on progressive transfer learning to incorporate knowledge gained from a sequence of datasets in a LSTM-based LM [10].

Let there be m number of corpora from which the knowledge is gained. A LSTM model built on corpus D_i to learn $p(\mathbf{x}|D_i)p(D_i)$ will have its parameters ω_i . It can be expressed as follows.

$$\begin{aligned} \prod_{i=1}^m p(\mathbf{x}|D_i)p(D_i) &\approx \prod_{i=1}^m p(\mathbf{x}|D_i, \omega_i)p(D_i, \omega_i) \\ &= \prod_{i=1}^m p(\mathbf{x}|D_i, \omega_i)p(\omega_i|D_i)p(D_i) \end{aligned} \quad (4)$$

If the same LM model is sequentially built from the given m datasets, parameters ω_i learned on i^{th} dataset will only depend on the parameters ω_{i-1} learned on the $(i-1)^{th}$ dataset, applying the Markov assumption.

$$\prod_{i=1}^m p(\mathbf{x}|D_i, \omega_i)p(\omega_i|D_i)p(D_i) \approx \prod_{i=1}^m p(\mathbf{x}|D_i, \omega_i)p(\omega_i|D_i, \omega_{i-1})p(D_i) \quad (5)$$

Here ω_0 is the initial weight that might be assigned randomly. Assuming the same probability (or uncertainty) for each dataset, transfer learning can be expressed as follows.

$$\begin{aligned} p(\mathbf{x}, D_1, \dots, D_n) &\approx \prod_{i=1}^m p(\mathbf{x}|D_i, \omega_i)p(\omega_i|D_i, \omega_{i-1})p(D_i) \\ &= \prod_{i=1}^m p(\mathbf{x}|D_i, \omega_i)p(\omega_i|D_i, \omega_{i-1}) \\ &\propto \sum_{i=1}^m \ln(p(\mathbf{x}|D_i, \omega_i)p(\omega_i|D_i, \omega_{i-1})) \end{aligned} \quad (6)$$

Followings can be inferred from Equation 6. (1) Each dataset D_i relevant to the application domain of LM can reduce uncertainty [10,11]. (2) Pr-training of LSTM-LM should be done by the order of the dataset of general population distribution to the dataset of specific population distribution because the parameter vector ω_i depends on ω_{i-1} [9,10]. For example, we can approximate the population distribution of Queensland (i.e., specific) from that of Australia (i.e., general) but the opposite is not true.

3.2 Regularising Classifier with Language Model

For the downstream task of classification, a LSTM is trained to learn $p(C_k, \mathbf{x}) \approx p(C_k | \mathbf{x}, \theta) p(\mathbf{x})$ on a small training dataset. The prior distribution $p(\mathbf{x})$ can be learned by sequentially pretraining an LSTM on m unlabelled datasets from general to specific domain. Using Equation 6, we can write the regularised classifier as follows.

$$\begin{aligned} p(C_k, \mathbf{x}) &\approx p(C_k, \mathbf{x}) p(\mathbf{x}, \omega) \\ &\approx p(C_k | \mathbf{x}, \theta) \prod_{i=1}^m p(\mathbf{x} | D_i, \omega_i) p(\omega_i | D_i, \omega_{i-1}) \end{aligned} \quad (7)$$

Figure 1 shows the process of transfer learning through LSTM-based LM to a LSTM classification model. Layers 1 to 3 are stacked LSTM layers. The LSTM-based LM (the left hand side model in Figure 1) is generated using three stacked layers along with embedding layer and LM softmax layer. The LM softmax layer is active during pretraining of LM with a sequence of unlabelled datasets and then it is frozen. Once the LM is pretrained, the class softmax layer and linear layer are augmented, as shown by the right hand side model in Figure 1. These two layers along with the pretrained LM active layers are trained with the small labelled dataset to learn the classification task of bias detection.

The main task of these additional two layers is to learn $p(C_k | \mathbf{x}, \theta)$. The combined network learns $p(C_k | \mathbf{x}, \theta) p(\mathbf{x}, \omega)$. The additional two layers are augmented at the end of NNLM in this model to assure that θ is learned from the finetuning of ω with labelled dataset during classification training. This process of training a classifier model (e.g. LSTM) generates a classifier regularised by the language model (LM) based transfer learning [10]. We call the trained model as LSTM-LM.

4 Empirical Analysis

Extensive experiments were conducted to evaluate the accuracy of the proposed method for bias detection. We used six standard classification evaluation measures [10]: Accuracy (Ac), Precision (Pr), Recall (Re), F₁ Score (F₁), Cohen Kappa Score (CKS) and Area Under Curve (AUC).

4.1 Data Collection

The iShield.ai Dataset (Version 1⁵) has a total data count of 57,424 sentences, where 27,131 sentences are biased (47%) and 30,293 sentences are unbiased (53%). There are four different kinds of biases: (1) 20,690 (36%) sentences

⁵ In the next versions, this dataset has been considerably modified.

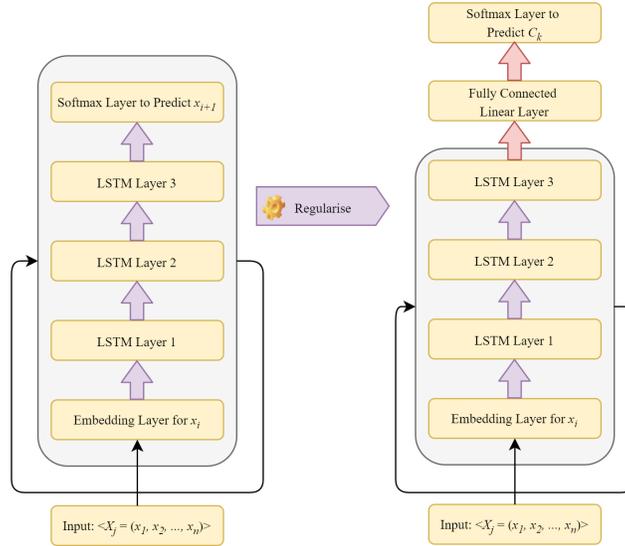


Fig. 1. LSTM-L: Process of regularising a classifier using a language model

are GENDER biased, (2) 4,339 (7.56%) sentences are RACE biased, (3) 1,553 (2.7%) sentences are ambiguous and (4) 549 (0.96%) sentences are AGE biased. We used 80-10-10% ratio between training, validation and testing instances. The dataset comes from two sub-domains namely Job Description (JD) and Non Job Description (NJD) that contributed 25,123 (43.75%) and 32,301 (56.25%) sentences respectively.

We used the following four definitions of bias in this research. **GENDER** – Any conscious or unconscious attempt to single out a particular gender or people who identify with a particular gender identity. Example: The claims assistant will handles the claims of veterans and their wives. **RACE** – Any conscious or unconscious attempt to single out a particular race or it’s people. Example: Own development of brownbag sessions and facilitate publishing reusable content to the organisation. **AGE** – Any conscious or unconscious attempt to single out a particular age bracket or it’s people. Example: We are a young organisation looking for young and talented marketers. **Ambiguous** – Any part of a sentence which does not clearly convey the intended meaning. Example: We are looking for a smart candidate for this position.

The basic unit of division used for annotation is a sentence. Any document is first split in sentences. Then sentences are passed to annotators in batches, where each batch consists of 150 sentences. Biases are usually observed as a group of four to five words in a sentence. A three member panel is set up for a Quality Assurance pass. The panel evaluates the labelled sentences for two checks. (1) The labelled sentences are adhering to the outlined definitions of bias. (2) Conflicting instances (e.g. a sentence should have been labelled as biased but is not) are eliminated or placed in the correct category.

Sexist Statement in Workplace (SSW) Dataset [23] This dataset was collected to check how effectively the classification model works on other datasets besides the bias detection dataset. The SSW dataset has around 1100 labelled instances for sexism statement in workplace. The instances are roughly balanced between sexism (labelled 1) and neutral (labelled 0) cases. Some examples from this dataset is given in Table 1. We used 80-20% ratio between training and testing instances. Hyper parameters were set by using cross validation in the 80% data used for training.

Table 1. Instances from SSW Dataset

Label	Statement
1	Women always get more upset than men.
1	The people at work are childish. it's run by women and when womendont agree to something, oh man.
1	I'm going to miss her resting bitch face.
0	No mountain is high enough for a girl to climb.
0	It seems the world is not ready for one of the most powerful andinfluential countries to have a woman leader. So sad.
0	Can you explain why what she described there is wrong?

Model Pretraining Datasets We collected a list of pretrained word vectors from [37]. The list has one million word vectors trained on Wikipedia 2017, UMBC webbase corpus and statmt.org news dataset (16B tokens). We use the following three corpora (D_1 , D_2 , and D_3) for sequentially pretraining LSTM-LM model and fineturning to target task, i.e. first LSTM-LM is pretrained using D_1 , then using D_2 and finally fintuned to target task using D_3 .

D_1 : The goal of using this corpus is to capture general properties of the English language. We pretrain the LSTM-LM model on Wikitext-103 that contains 28,595 preprocessed Wikipedia articles and 103 million words [36]. After pretraining on D_1 , we approximate the probability distribution $p(\mathbf{x}|D_1, \omega_1)$.

D_2 : The goal of using this corpus is to bridge the data distribution between the target task domain (i.e., bias detection in job description) and the general domain (i.e. standard language). This is because the target task is likely to come from a different distribution than the general corpus. D_2 should be chosen such that it has commonalities with both D_1 reflecting a general domain (Wikipedia) and the corpus D_3 reflecting a target domain (e.g. labelled data of job description). We use a set of unlabelled JD and NJD data as D_2 .

4.2 Baseline Models

We have implemented 10 baseline models to compare the performance of the proposed LSTM-LM.

Models with pretrained word vectors (i.e. word embeddings by word2Vec) include (1) LSTM with pretrained Word vectors (LSTM-W) [26] and (2) CNN with pretrained Word vectors (CNN-W) [11]. LSTM-W has 100 units, 50%

dropout, binary cross-entropy loss function, Adam optimiser and sigmoid activation. The hyper parameters of CNN-W is set as in [11]. We used one million pretrained word vectors each with 300-dimension [37]. Word vectors are pretrained on Wikipedia 2017, UMBC webbase corpus and statmt.org news dataset. A Continuous Bag-of-Words Word2vec [38] model is used in pretraining.

LSTM and CNN without pretrained word vectors (LSTM-P) [26,11]. These are traditional LSTM and CNN models that have not been pretrained by any data. Similar to LSTM-W, LSTM-P has 100 units, 50% dropout, binary cross-entropy loss function, Adam optimiser and sigmoid activation. The hyper parameters of CNN is set as in [11].

Feedforward Deep Neural Network (DNN) [22]. It has five hidden layers, each layer containing eighty units, 50% dropout applied to the input layer and the first two hidden layers, softmax activation and 0.04 learning rate. For all neural network based models, hyperparameters are manually tuned based on cross-validation.

Non NN models include Support Vector Machines (SVM) [24], Random Forest (RF) [33], Decision Tree (DT) [44], Gaussian Naive Bayes (GNB) [32], k-Nearest Neighbours (kNN) [46] and Ridge Classifier (RC) [27]. Hyperparameters of all these models are automatically tuned using ten-fold cross-validation and GridSearch using sklearn library.

None of the models, except LSTM-LM, LSTM-W and CNN-W, are pretrained or utilised unlabelled dataset.

4.3 Experimental Results: SSW Dataset

SSW is a very small dataset. The experimental results on SSW dataset are given in Table 2. The proposed language model-based transfer learning model LSTM-LM outperforms all the baseline models. Beside LSTM-LM, other two word vector-based transfer learning-based models, CNN-W and LSTM-W yield the second and third best performance, respectively. We conjecture that (1) transfer learning-based models produce better outcome, and (2) the language model-based transfer learning brings more benefits than the word vector-based transfer learning when the training dataset is small.

It is interesting to note that CNN-W outperforms CNN as well as LSTM-W outperforms LSTM. CNN-W and CNN (or LSTM-W and LSTM) use the exactly same architecture, except that CNN-W (or LSTM-W) uses pretrained word vectors for transfer learning. This further emphasises the benefit of using transfer learning over the standard models when the training dataset is small.

Traditional models (i.e. RF, DT, GNB and kNN) do not utilise any transfer learning and solely rely on the labelled training dataset. Therefore, they give lower performance when the labelled training dataset is small.

4.4 Experimental Results: The iShield.ai Dataset

The experimental results comparing our LSTM-LM against ten baseline models on the iShield dataset are given in Table 3.

Table 2. Experimental Results on SSW Dataset

	Sample Average			Weighted Average			Macro Average			
	Accuracy	AUC	CKS	Precision	Recall	F ₁ -score	Precision	Recall	F ₁ -score	Support
LSTM-LM	0.89	0.89	0.77	0.89	0.89	0.89	0.89	0.88	0.89	228
RF	0.82	0.81	0.63	0.82	0.82	0.82	0.81	0.81	0.81	228
DT	0.79	0.78	0.57	0.79	0.79	0.79	0.78	0.78	0.78	228
GNB	0.70	0.71	0.40	0.72	0.7	0.7	0.71	0.71	0.7	228
kNN	0.58	0.60	0.19	0.63	0.58	0.56	0.62	0.6	0.57	228
SVM	0.82	0.81	0.63	0.82	0.82	0.82	0.82	0.81	0.81	228
RC	0.77	0.77	0.53	0.77	0.77	0.77	0.77	0.77	0.77	228
CNN-W	0.86	0.86	0.72	0.86	0.86	0.86	0.86	0.86	0.86	228
CNN	0.82	0.81	0.63	0.82	0.82	0.82	0.82	0.81	0.82	228
LSTM-W	0.85	0.84	0.70	0.85	0.85	0.85	0.86	0.84	0.85	228
LSTM	0.82	0.82	0.64	0.82	0.82	0.82	0.82	0.82	0.82	228

Accuracy, AUC and CKS are three important measures for understanding the overall significance of a classification model. Table 3 shows that LSTM-LM gives the best Accuracy, AUC and CKS results. CKS indicates the reliability between the prediction made by a model and the ground truth. All the three transfer learning-based models (i.e. LSTM-LM, LSTM-W and CNN-W) have high CKS value. However, LSTM-LM has better accuracy and AUC than other two (i.e. LSTM-W and CNN-W).

Overall LSTM-LM gives us the best results, as indicated by best weighted average precision, recall and F₁ score. Even though SVM, CNN-W and LSTM-W have the same weighted average precision value as LSTM-W, their recall and F₁ score are lower than LSTM-W.

High precision indicates most of the identified biased sentences are indeed bias. However, if the recall is not high enough, then many bias sentences will left undetected. Therefore, better recall is desirable in bias detection. Yet the excessive false positives can result in a higher cost for investigating many false detection. Therefore, a balance in both recall and precision is needed. A higher F₁ score indicates both precision and recall is high. LSTM-LM gives the best F₁ score for both weighted average and macro average. Macro Average is used to evaluate the performance of a classifier for minority classes (classes with small number of instances), where weighted average favours majority classes. The high F₁-score for both weighted average and macro average indicates that LSTM-LM works reasonably well for both majority and minority classes. Best macro average precision is achieved by SVM, but SVM has very poor recall value and F₁-score.

5 Deployment and Architecture of Integrated System

The proposed model LSTM-LM is deployed in the flavour of following four different products for the convenience of end users by iShield.ai. (1) **Dost**⁶: This bot can be configured to detect bias on enterprise communication platforms such as *Slack* and *Microsoft teams*. (2) **Chrome Plugin**⁷: This plugin can be configured

⁶ <https://ishield.ai/dost>

⁷ <https://ishield.ai/chrome-plugin>

Table 3. Experimental Results on the iShield.ai dataset

	Sample Average			Weighted Average			Macro Average			
	Accuracy	AUC	CKS	Precision	Recall	F ₁ -score	Precision	Recall	F ₁ -score	Support
LSTM-LM	0.85	0.78	0.73	0.84	0.85	0.84	0.69	0.6	0.63	5743
RF	0.84	0.76	0.72	0.83	0.84	0.83	0.72	0.57	0.61	5743
DT	0.83	0.77	0.71	0.83	0.83	0.83	0.64	0.59	0.61	5743
GNB	0.62	0.74	0.45	0.79	0.62	0.67	0.44	0.58	0.43	5743
kNN	0.79	0.70	0.64	0.78	0.79	0.78	0.66	0.48	0.52	5743
SVM	0.84	0.73	0.72	0.84	0.84	0.83	0.81	0.51	0.55	5743
RC	0.83	0.75	0.70	0.82	0.83	0.82	0.68	0.57	0.6	5743
CNN-W	0.84	0.76	0.73	0.84	0.84	0.83	0.73	0.57	0.59	5743
CNN	0.84	0.75	0.72	0.83	0.84	0.83	0.58	0.56	0.56	5743
LSTM-W	0.84	0.77	0.73	0.84	0.84	0.84	0.67	0.58	0.6	5743
LSTM	0.84	0.76	0.72	0.84	0.84	0.83	0.69	0.58	0.6	5743

with chrome browser for screening text contents for bias in any web applications.

(3) **Content screener**⁸: This web application allows Checking for bias in created text contents before publishing them. (4) **Application Programming Interface (API)**⁹: This API can be integrated with enterprise platforms where contents are created and published.

The backend of the system is deployed on Amazon Web Services (AWS) and is built to scale for industry use. The current architecture can process 8 parallel requests. This architecture can be easily scaled to accommodate increased request volumes. The architecture of the backend system is shown in Fig. 2 and can be explained as follows. (1) Multiple user requests are received concurrently. (2) Requests are received at the AWS Lambda Function. (3) Each model is placed in an n-core EC2 instance. (3.1) Gunicorn Web Server Gateway Interface is implemented with each model to gather and distribute requests for parallel processing. At the time of writing this paper, 8 parallel requests can be processed by each Gunicorn WSGI. This can be scaled up with increased volumes. (3.2) Each model is placed in an independent docker container to isolate it’s function from other environment related dependencies. (4) Once a piece of content is identified *biased* by a model, asynchronous requests travel back for confidence based sorting. (5) After index calculation, a response JSON is prepared and sent to the AWS Lamba function. (6) Database operations are performed at the AWS Lamda Function, after which the results travel back to the user.

6 Conclusion

We propose a transfer learning based language model to learn from manually tagged documents for automatically identifying bias in enterprise content in order to create the workplace more inclusive. The trained model is thoroughly evaluated on independent datasets to ensure a general application, and it is deployed in a real-world application.

⁸ <https://ishield.ai/screener>

⁹ <https://ishield.ai/api>

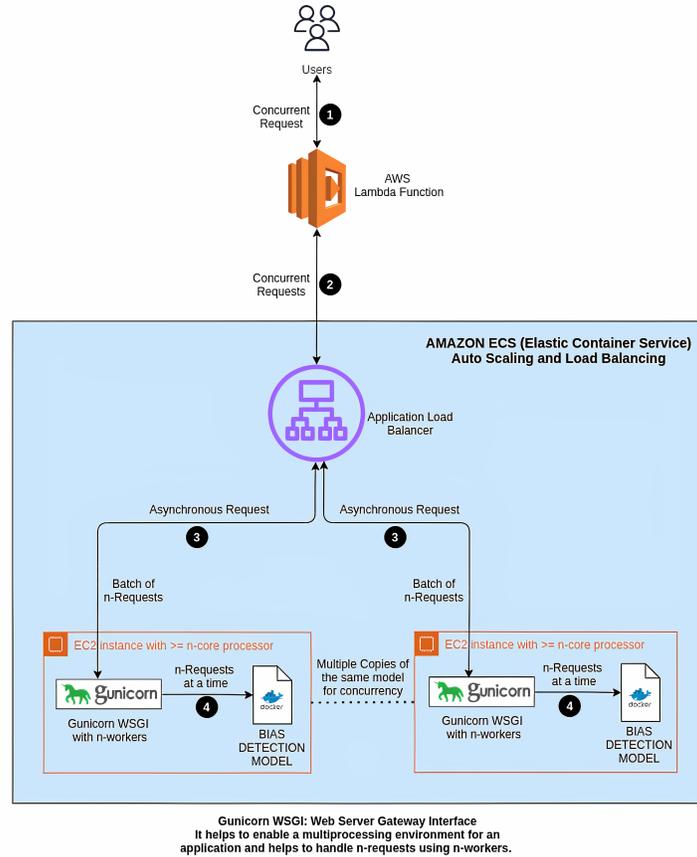


Fig. 2. Bias Detection System Architecture

References

1. The Psychology of Inclusion and the Effects in Advertising: Gen Z. Tech. rep., Microsoft Advertising (2020)
2. Abul, B., NayakRichi: Active Learning for Effectively Fine-Tuning Transfer Learning to Downstream Task. ACM Transactions on Intelligent Systems and Technology (TIST) **12**(2) (2 2021). <https://doi.org/10.1145/3446343>, <https://dl.acm.org/doi/abs/10.1145/3446343>
3. Agovino, T.: Toxic Workplace Cultures Are Costing Employers Billions. Tech. rep., Society for Human Resource Management (2019), <https://www.talkworkculture.com/advice-info/toxic-workplace/>
4. Attenberg, J., Provost, F.: Why label when you can search?: alternatives to active learning for applying human resources to build classification models under extreme class imbalance. In: Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining. pp. 423–432. ACM (2010)

5. Attenberg, J., Provost, F.: Inactive learning?: difficulties employing active learning in practice. *ACM SIGKDD Explorations Newsletter* **12**(2), 36–41 (2011)
6. Bailinson, P., Decherd, W., Ellsworth, D., Guttman, M.: Understanding organizational barriers to a more inclusive workplace. Tech. rep., McKinsey Insights (2020)
7. Bashar, M., Nayak, R.: QutNocturnal@HASOC’19: CNN for hate speech and offensive content identification in Hindi language. In: *CEUR Workshop Proceedings*. vol. 2517 (2019)
8. Bashar, M., Nayak, R., Suzor, N., Weir, B.: Misogynistic tweet detection: Modelling CNN with small datasets, vol. 996 (2019). https://doi.org/10.1007/978-981-13-6661-1_11
9. Bashar, M.A., Nayak, R., Luong, K., Balasubramaniam, T.: Progressive domain adaptation for detecting hate speech on social media with small training set and its application to COVID-19 concerned posts. *Social Network Analysis and Mining* 2021 11:1 **11**(1), 1–18 (7 2021). <https://doi.org/10.1007/S13278-021-00780-W>, <https://link.springer.com/article/10.1007/s13278-021-00780-w>
10. Bashar, M.A., Nayak, R., Suzor, N.: Regularising LSTM classifier by transfer learning for detecting misogynistic tweets with small training set. *Knowledge and Information Systems* pp. 1–26 (6 2020). <https://doi.org/10.1007/s10115-020-01481-0>, <https://link.springer.com/article/10.1007/s10115-020-01481-0>
11. Bashar, M.A., Nayak, R., Suzor, N., Weir, B.: Misogynistic Tweet Detection: Modelling CNN with Small Datasets. In: *The 16th Australasian Data Mining Conference*. pp. 3–16. Springer (2018)
12. Bordia, S., Bowman, S.R.: Identifying and Reducing Gender Bias in Word-Level Language Models. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, Stroudsburg, PA, USA (2019). <https://doi.org/10.18653/v1/N19-3002>, <http://aclweb.org/anthology/N19-3002>
13. Bradbury, J., Merity, S., Xiong, C., Socher, R.: Quasi-recurrent neural networks. *arXiv preprint arXiv:1611.01576* (2016)
14. C Lin, M.: Active Learning with Unbalanced Classes & Example-Generated Queries. In: *AAAI Conference on Human Computation* (2018)
15. Craven, M., Kumlien, J., others: Constructing biological knowledge bases by extracting information from text sources. In: *ISMB*. vol. 1999, pp. 77–86 (1999)
16. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018)
17. Dinan, E., Fan, A., Williams, A., Urbanek, J., Kiela, D., Weston, J.: Queens are Powerful too: Mitigating Gender Bias in Dialogue Generation. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. pp. 8173–8188. Association for Computational Linguistics, Stroudsburg, PA, USA (11 2020). <https://doi.org/10.18653/v1/2020.emnlp-main.656>, <https://www.aclweb.org/anthology/2020.emnlp-main.656>
18. Dinan, E., Fan, A., Wu, L., Weston, J., Kiela, D., Williams, A.: Multi-Dimensional Gender Bias Classification. Tech. rep.
19. Dixon-Fyle, S., Dolan, K., Hunt, V., Prince, S.: Diversity wins: How inclusion matters. Tech. rep., McKinsey & Company (2020), <https://www.mckinsey.com/featured-insights/diversity-and-inclusion/diversity-wins-how-inclusion-matters>
20. Ehrenberg, H.R., Shin, J., Ratner, A.J., Fries, J.A., Ré, C.: Data programming with ddlite: Putting humans in a different part of the loop. In: *Proceedings of the Workshop on Human-In-the-Loop Data Analytics*. p. 13. ACM (2016)

21. Garr, S.S., Jackson, C.: Diversity & inclusion technology: The rise of a transformative market. Tech. rep., Mercer, New York, United States (2019), <https://www.mercer.com/our-thinking/career/diversity-and-inclusion-technology.html>
22. Glorot, X., Bengio, Y.: Understanding the difficulty of training deep feedforward neural networks. In: Proceedings of the thirteenth international conference on artificial intelligence and statistics. pp. 249–256 (2010)
23. Grosz, D., Conde-Cespedes, P.: Automatic Detection of Sexist Statements Commonly Used at the Workplace. Tech. rep. (2020), <https://hal.archives-ouvertes.fr/hal-02573576>
24. Hearst, M.A., Dumais, S.T., Osuna, E., Platt, J., Scholkopf, B.: Support vector machines. *IEEE Intelligent Systems and their applications* **13**(4), 18–28 (1998)
25. Herbert, F.: Is unconscious bias training still worthwhile? *LSE Business Review* (3 2021), <https://blogs.lse.ac.uk/businessreview/>
26. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural computation* **9**(8), 1735–1780 (1997)
27. Hoerl, A.E., Kennard, R.W.: Ridge regression: applications to nonorthogonal problems. *Technometrics* **12**(1), 69–82 (1970)
28. Hube, C., Fetahu, B.: Detecting Biased Statements in Wikipedia. In: The Web Conference 2018 - Companion of the World Wide Web Conference, WWW 2018. pp. 1779–1786. Association for Computing Machinery, Inc, New York, New York, USA (4 2018). <https://doi.org/10.1145/3184558.3191640>, <http://dl.acm.org/citation.cfm?doid=3184558.3191640>
29. Johnson, S.K., Hekman, D.R., Chan, E.T.: If There’s Only One Woman in Your Candidate Pool, There’s Statistically No Chance She’ll Be Hired. Tech. rep., Harvard Business Review, <https://hbr.org/2016/04/if-theres-only-one-woman-in-your-candidate-pool-theres-statistically-no-chance-shell-be-hired>
30. Jozefowicz, R., Vinyals, O., Schuster, M., Shazeer, N., Wu, Y.: Exploring the limits of language modeling. arXiv preprint arXiv:1602.02410 (2016)
31. Kaneko, M., Bollegala, D.: Gender-preserving Debiasing for Pre-trained Word Embeddings. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. pp. 1641–1650. Association for Computational Linguistics, Stroudsburg, PA, USA (2019). <https://doi.org/10.18653/v1/P19-1160>, <https://www.aclweb.org/anthology/P19-1160>
32. Lewis, D.D.: Naive (Bayes) at forty: The independence assumption in information retrieval. In: European conference on machine learning. pp. 4–15. Springer (1998)
33. Liaw, A., Wiener, M., others: Classification and regression by randomForest. *R news* **2**(3), 18–22 (2002)
34. Manzini, T., Yao Chong, L., Black, A.W., Tsvetkov, Y.: Black is to Criminal as Caucasian is to Police: Detecting and Removing Multiclass Bias in Word Embeddings. In: Proceedings of the 2019 Conference of the North. vol. 1, pp. 615–621. Association for Computational Linguistics, Stroudsburg, PA, USA (2019). <https://doi.org/10.18653/v1/N19-1062>, <http://aclweb.org/anthology/N19-1062>
35. Maudslay, R.H., Gonen, H., Cotterell, R., Teufel, S.: It’s All in the Name: Mitigating Gender Bias with Name-Based Counterfactual Data Substitution. EMNLP-IJCNLP 2019 - 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing, Proceedings of the Conference pp. 5267–5275 (9 2019), <http://arxiv.org/abs/1909.00871>

36. Merity, S., Xiong, C., Bradbury, J., Socher, R.: Pointer sentinel mixture models. arXiv preprint arXiv:1609.07843 (2016)
37. Mikolov, T., Grave, E., Bojanowski, P., Puhersch, C., Joulin, A.: Advances in Pre-Training Distributed Word Representations. In: Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018) (2018)
38. Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. In: Advances in neural information processing systems. pp. 3111–3119 (2013)
39. Mikolov, T., Karafiát, M., Burget, L., Černocký, J., Khudanpur, S.: Recurrent neural network based language model. In: Eleventh Annual Conference of the International Speech Communication Association (2010)
40. Mitchell, M., Research, G., Baker, D., Denton, E., Hutchinson, B., Hanna, A., Morgenstern, J., Moorosi, N., Gebru, T.: Diversity and Inclusion Metrics in Subset Selection. Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society **7** (2020). <https://doi.org/10.1145/3375627>, <https://doi.org/10.1145/3375627.3375832>
41. Patterson, G., Van Horn, G., Belongie, S.J., Perona, P., Hays, J.: Tropel: Crowdsourcing Detectors with Minimal Training. In: HCOMP. pp. 150–159 (2015)
42. Pryzant, R., Martinez, R.D., Dass, N., Kurohashi, S., Jurafsky, D., Yang, D.: Automatically neutralizing subjective bias in text. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 480–489 (4 2020). <https://doi.org/10.1609/aaai.v34i01.5385>, <https://github.com/rpryzant/neutralizing-bias>
43. Recasens, M., Danescu-Niculescu-Mizil, C., Jurafsky, D.: Linguistic Models for Analyzing and Detecting Biased Language. In: Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics. pp. 1650–1659. Association for Computational Linguistics, Sofia, Bulgaria (2013), <http://en.wikipedia.org/wiki/>
44. Safavian, S.R., Landgrebe, D.: A Survey of Decision Tree Classifier Methodology. IEEE Transactions on Systems, Man and Cybernetics **21**(3), 660–674 (1991). <https://doi.org/10.1109/21.97458>
45. Wang, T., Zhao, J., Yatskar, M., Chang, K.W., Ordonez, V.: Balanced Datasets Are Not Enough: Estimating and Mitigating Gender Bias in Deep Image Representations. Tech. rep.
46. Weinberger, K.Q., Saul, L.K.: Distance metric learning for large margin nearest neighbor classification. Journal of Machine Learning Research **10**(Feb), 207–244 (2009)
47. Zalis, S.: Inclusive ads are affecting consumer behavior, according to new research. Tech. rep., Think with Google (2019), <https://www.thinkwithgoogle.com/future-of-marketing/management-and-culture/diversity-and-inclusion/thought-leadership-marketing-diversity-inclusion/>
48. Zmigrod, R., Mielke, S.J., Wallach, H., Cotterell, R.: Counterfactual Data Augmentation for Mitigating Gender Stereotypes in Languages with Rich Morphology. ACL 2019 - 57th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference pp. 1651–1661 (6 2019), <http://arxiv.org/abs/1906.04571>